

Coarse-Grained Topology Estimation via Graph Sampling*

Maciej Kurant
CalIT2
UC Irvine
mkurant@uci.edu

Zack W. Almquist
Sociology Dept, CalIT2
UC Irvine
almquist@uci.edu

Minas Gjoka
CalIT2
UC Irvine
mgjoka@uci.edu

Carter T. Butts
Sociology Dept, CalIT2, IMBS
UC Irvine
buttsc@uci.edu

Yan Wang
CalIT2
UC Irvine
wang.yan@uci.edu

Athina Markopoulou
EECS, CalIT2, CPCC
UC Irvine
athina@uci.edu

ABSTRACT

Many online networks are measured and studied via sampling techniques, which typically collect a relatively small fraction of nodes and their associated edges. Past work in this area has primarily focused on obtaining a representative sample of nodes and on efficient estimation of local graph properties (such as node degree distribution or any node attribute) based on that sample. However, less is known about estimating the global topology of the underlying graph.

In this paper, we show how to efficiently estimate the coarse-grained topology of a graph from a probability sample of nodes. In particular, we consider that nodes are partitioned into *categories* (e.g., countries or work/study places in OSNs), which naturally defines a weighted *category graph*. We are interested in estimating (i) the size of categories and (ii) the probability that nodes from two different categories are connected. For each of the above, we develop a family of estimators for design-based inference under uniform or non-uniform sampling, employing either of two measurement strategies: *induced subgraph sampling*, which relies only on information about the sampled nodes; and *star sampling*, which also exploits category information about the neighbors of sampled nodes. We prove consistency of these estimators and evaluate their efficiency via simulation on fully known graphs. We also apply our methodology to a sample of Facebook users to obtain a number of category graphs, such as the college friendship graph and the country friendship graph; we share and visualize the resulting data at www.geosocialmap.com.

Keywords

Online Social Networks, coarse-grained topology, induced subgraph sampling, star sampling, Facebook.

1. INTRODUCTION

Many large online networks, such as online social networks (OSNs) and the World Wide Web (WWW), are

*We make our datasets available, together with a customizable web-based visualization at www.geosocialmap.com

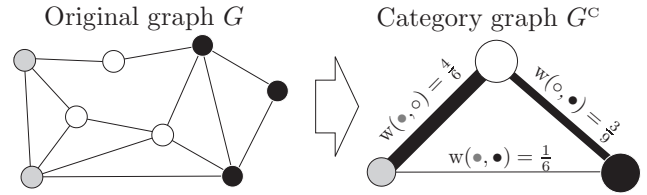


Figure 1: Nodes in the original graph (G) belong in one of three categories: white, gray, and black. The category graph (G^c) consists of three nodes, corresponding to the three categories, connected by weighted edges. The edge weight $w(o, \bullet)$ in G^c is the probability that a black and a white node, randomly chosen from G , are connected in G (see Eq.(3)). The main goal of this paper is to estimate these edge weights based on a probability sample of nodes of G .

currently studied via sampling techniques. Sampling becomes necessary due to the sheer size of these networks and/or access limitations, which make it infeasible to collect (and, in some cases, to analyze) these networks in their entirety.

Most principled graph sampling methods to date have focused on collecting a probability sample of nodes [6, 19, 20, 30, 35, 51–53, 60]. Based on such a sample, one can efficiently estimate many local graph properties, such as node attribute frequency, degree distribution, degree-degree correlations, or clustering coefficients [26, 34]. However, these features reveal little about the global properties of the underlying graph, such as path-based properties (connectivity, diameter, average shortest path length) or community structure.

In this paper, we show how a particular aspect of global network structure, namely coarse-grained topology, can be efficiently estimated from a probability sample of nodes. Specifically, we note that nodes in many online graphs belong to *categories*, explicitly declared by users or clearly determined by observable characteristics. For example, in Facebook, users can officially declare the college or workplace with which they are af-

filiated, or a country/city in which they live. Similarly, in the WWW, all nodes can be categorized by their domain names, and the users of Internet radio sites like Last.FM may be grouped on the basis of listening behavior. This potentially allows us to build and study *category graphs*, in which each node corresponds to a category and edge weights reflect the frequency of ties between category members in the original graph. We illustrate these concepts in Fig. 1.

The contribution of this paper lies in developing and evaluating several efficient estimators for two properties of the category graph, namely the size of the categories and the edge weights. These estimators take as input a uniform or non-uniform probability sample of nodes, measured via one of two strategies: *induced subgraph sampling*, in which we have information regarding only the sampled nodes; and *star sampling*, in which we also have category information about the neighbors of sampled nodes. We show that our estimators have good asymptotic properties (consistency, and hence asymptotic unbiasedness) and we evaluate their efficiency via simulation: employing fully observed graphs from both synthetic and empirical sources, we examine how estimator performance varies with the properties of the underlying graph. Finally, as a practical illustration of our approach, we apply our methodology to a sample of Facebook nodes to estimate several Facebook category graphs, such as the inter-college and inter-country friendship graphs. The resulting Facebook category graphs are made available (along with a highly-customizable, web-based visualization service) at www.geosocialmap.com.

The structure of the remainder of the paper is as follows. Section 2 presents the problem statement. Section 3 reviews node sampling techniques. Sections 4 and 5 present our estimators for uniform and non-uniform probability samples, respectively. Section 6 presents simulation results on fully known graphs. Section 7 applies our estimators to samples of Facebook. Section 8 reviews related work. Section 9 concludes the paper. Finally, in Appendix we prove the consistency of all estimators proposed in this paper.

2. NOTATION AND PROBLEM STATEMENT

2.1 Basic graph G

We consider an undirected, static¹ graph $G = (V, E)$, with $N = |V|$ nodes and $|E|$ edges. Denote by $\deg(v)$

¹Sampling dynamic graphs is currently an active research area [51,60,67], but out of the scope of this paper. Indeed, during the collection of Facebook data sets we use, the underlying graphs changed very insignificantly [20,35]. Moreover, in this paper we focus on coarse granularity, which should change even more slowly in time, as argued in [67].

the degree of node $v \in V$, and by

$$\text{vol}(A) = \sum_{v \in A} \deg(v) \quad (1)$$

the volume of a set of nodes $A \subseteq V$. We will often use

$$f_A = \frac{|A|}{|V|} \quad \text{and} \quad f_A^{\text{vol}} = \frac{\text{vol}(A)}{\text{vol}(V)} \quad (2)$$

to denote the relative size of A in terms of number of nodes and volume, respectively.

2.2 Category graph G^c

We assume that the set of nodes V is partitioned into a set \mathcal{C} of *categories*, i.e., that $\bigcup_{C \in \mathcal{C}} C = V$. We are interested in the *category graph* $G^c = (\mathcal{C}, E^c)$, with node set given by the categories of G .² For two different categories $A, B \in \mathcal{C}$, $A \neq B$, denote by $E_{A,B} \subset E$ the corresponding edge-cut in G , i.e.,

$$E_{A,B} = \{\{u, v\} \in E : u \in A \text{ and } v \in B\}.$$

If $|E_{A,B}| > 0$ then we draw an edge $\{A, B\}$ between A and B in G^c . We show an example of a category graph in Fig. 1.

The way we defined category graph G^c so far, prevents self-loops, but potentially allows for edge weights. The *weight* $w(A, B)$ of edge $\{A, B\}$ can be defined in a number of ways. For instance, one could trivially set it always equal to 1. In some settings, e.g., statistical modeling, the number of inter-category edges, $w(A, B) = |E_{A,B}|$, is a useful choice. For many purposes, however, it is useful to have a notion of edge weight that adjusts for category size, e.g.,

$$w(A, B) = \frac{|E_{A,B}|}{|A| \cdot |B|}. \quad (3)$$

This definition has an intuitive interpretation. Because $|A| \cdot |B|$ is the size of the maximum possible edge-cut from A to B , $w(A, B)$ is equal to the probability that a uniformly selected member of A is connected to a uniformly selected member of B . We give an example of these weights $w(A, B)$ in Fig. 1.

2.3 Goal: Estimate G^c through sampling

Given the full knowledge of graph G , it is trivial to construct the category graph with its edge weights. In many cases, however, the knowledge of the full graph G is not available, rendering exact computation of Eq.(3) infeasible. For instance, downloading the entire Facebook social graph via HTML scraping would require

²We are not the first ones to be interested in coarse-grained structures. See, e.g., the social network literature on *block-models* [66], in which our categories correspond to *positions*, our category graph to the *reduced graph* or *block image*, and our edge weights to *block densities* or *mixing rates*. See Section 8 for additional references.

downloading and processing about 50 terabytes of HTML traffic [20], which is rather prohibitive in practice.

In contrast, it is often possible to collect a *sample* $S \subseteq V$ of nodes of G . Note that we permit S to contain multiple copies of the same node, *i.e.*, the sampling with replacement. The challenge, then, and the main goal of this paper is to estimate the category graph G^c based on the sample S .

3. SAMPLING

Our methodology takes as input a probability sample of nodes. Obtaining such a sample is an active research topic in its own right (see Section 8). In Section 3.1, we briefly review the node sampling techniques that we use later in simulations and Facebook implementation.

Independently of the sampling technique employed, we may collect less or more category information on each sampled node. In Section 3.2, we describe two scenarios most common in practice. As we will see later, they result in two different sets of estimators, often with very different performance.

3.1 Node sampling techniques

3.1.1 Independence Sampling

Under independence sampling, we sample nodes independently from the set V , with replacement. We distinguish two general cases: *Uniform Independence Sampling (UIS)*, where sampling probabilities are uniform (the same for all nodes); and *Weighted Independence Sampling (WIS)*, which samples v with probability proportional to a known weight $w(v)$.

In general, UIS and WIS are not feasible in online networks because of the lack of sampling frame. For example, the list of all user IDs may not be publicly available, or the user ID space may be too sparsely allocated to permit rejection sampling. Nevertheless, these techniques can occasionally be employed, either when permitted by fortuitous circumstances (see *e.g.*, use by [19,20]) or when deliberately “down-sampling” a large graph to speed analysis. Independence samplers are also conceptually important as a baseline for comparison with crawling-based sampling methods.

3.1.2 Sampling via Crawling

In contrast to independence sampling, crawling techniques are feasible in many online networks, and are thus the main focus of this paper. The crawling methods described here lead to an approximate probability sample (asymptotically approaching UIS or WIS) from the node set, in the limit of increasing sample size.

Simple Random Walk (RW) [41] selects the next-hop node v uniformly at random among the neighbors of the current node u . On a connected and aperiodic graph, RW samples node v with probability linearly proportional to its degree $\deg(v)$.

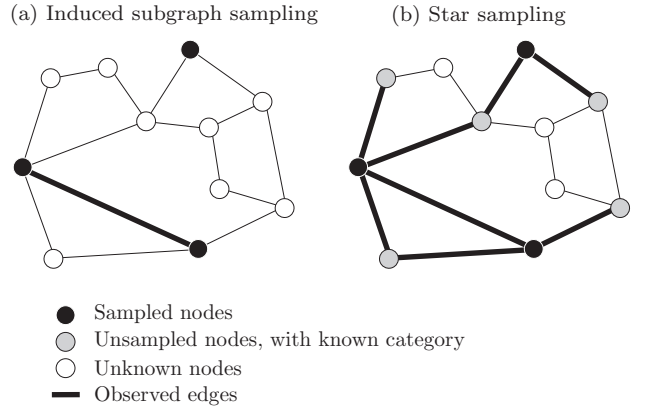


Figure 2: Observed categories and edges, under two scenarios we study in this paper.

Weighted Random Walk (WRW) is RW on a weighted graph [5]. In our simulations and implementation, we use “Stratified WRW,” or S-WRW [35], *i.e.*, a version of WRW that increases the sampling efficiency by over-sampling graph regions relevant to the measurement objective and under-sampling the irrelevant ones.

Metropolis-Hastings Random Walk (MHRW) is a version of random walk that modifies the transition probabilities to converge to a desired stationary distribution (often uniform). It was shown in [20,51] that RW outperforms MHRW for most applications, which we observe in our implementation as well.

3.2 Observed categories and edges

Our estimators will make use of every *fully observed edge*, *i.e.*, edge $\{u, v\}$ for which we know the categories of both u and v . We distinguish between two measurement scenarios [34] that yield different sets of observed edges, as follows.

3.2.1 Induced Subgraph Sampling

Under *induced subgraph sampling*, we learn the categories of the sampled nodes only. Consequently, the observed edges are only the edges induced on the set S of sampled nodes, as shown in Fig. 2(a).

3.2.2 Star Sampling

In some settings, sampling a node $u \in S$ reveals the categories of *all* its neighbors (not only the neighbors in S). This is typically the case when sampling is done through scraping the HTML pages of OSNs [20,35]. We refer to this as *star sampling*³ and we show an example in Fig. 2(b).

Finally, we emphasize that star sampling requires only information about neighbors’ *categories*; their degree or friend list is not needed, nor ties among neighbors (as in complete egonet sampling [66]).

³To be precise, following the terminology of [34], *labeled* star sampling. The *unlabeled* star sampling gets only the total number of neighbors, without their identities or categories.

4. UNIFORM SAMPLING

In this section, we provide design-based estimators for category sizes and category graph edge weights, given a uniform independence (UIS) sample from the node set. All estimators shown in this section and in Section 5 are consistent; proofs are provided in the Appendix.

4.1 Estimating category size ($|A|$)

Learning the size of a given category can be an important measurement objective per se. Moreover, it is also a building block of the edge weight estimators we derive in Section 4.2.2.

4.1.1 Induced subgraph sampling

The size $|A|$ of category A can be trivially estimated by multiplying by N the fraction of nodes sampled in A , *i.e.*,

$$|\hat{A}| = N \cdot \frac{|S_A|}{|S|}, \quad (4)$$

where

$$S_A = \{v \in S : v \in A\}$$

is a multiset containing all samples from category A .

4.1.2 Star sampling

Although not obvious at first blush, star sampling gives us an alternative way to estimate category sizes. Denote by

$$k_A = \frac{1}{|A|} \sum_{v \in A} \deg(v) \quad \text{and} \quad k_V = \frac{1}{|V|} \sum_{v \in V} \deg(v)$$

the average node degree in category A and in the entire graph, G , respectively. Because $\text{vol}(A) = |A| \cdot k_A$, we can re-write the relative volume f_A^{vol} of category A (see Eq.(2)) as

$$f_A^{\text{vol}} = \frac{\text{vol}(A)}{\text{vol}(V)} = \frac{|A| \cdot k_A}{|V| \cdot k_V} = \frac{|A| \cdot k_A}{N \cdot k_V}.$$

This allows us to estimate the size $|A|$ of category A as

$$|\hat{A}| = N \cdot \hat{f}_A^{\text{vol}} \cdot \frac{\hat{k}_V}{\hat{k}_A}. \quad (5)$$

This formula may seem less attractive than Eq.(4), because we now have to estimate three different numbers. However, k_V and k_A can be easily estimated, respectively by

$$\hat{k}_V = \frac{\sum_{v \in S} \deg(v)}{|S|} \quad \text{and} \quad \hat{k}_A = \frac{\sum_{v \in S_A} \deg(v)}{|S_A|}. \quad (6)$$

Similarly, f_A^{vol} could be estimated by

$$\hat{f}_A^{\text{vol}} = \frac{\sum_{v \in S} \deg(v) \cdot 1_{\{v \in A\}}}{\sum_{v \in S} \deg(v)}.$$

But we have proposed in [35] a much more efficient star-based estimator of f_A^{vol} , *i.e.*,

$$\hat{f}_A^{\text{vol}} = \frac{1}{\text{vol}(S)} \sum_{s \in S} \sum_{v \in \mathcal{N}(s)} 1_{\{v \in A\}}. \quad (7)$$

By plugging Eq.(6) and Eq.(7) into Eq.(5), we obtain a complex yet powerful star-based estimator of size $|A|$.

We show later that the star sampling estimator of Eq.(5) often outperforms the trivial estimator or Eq.(4), especially in dense graphs. One reason for this result is that Eq.(4) employs only the number $|S_A|$ of samples from A . This number is a random variable with a potentially high variance (especially for walks). In contrast, Eq.(5) relies on mean degree estimates rather than on counting-based estimates, which employ more information (edges not in $G[S]$) and tend to be more stable.

4.2 Estimating category edge weights ($w(A, B)$)

Recall from Eq.(3) that, given the full knowledge of graph G , the weight $w(A, B)$ is obtained by dividing the number of edges between A and B by the maximal possible number of such edges. We use this same idea when estimating $w(A, B)$ from our sample S , except that now we divide the number of edges *observed* between A and B by the maximal number of such edges we could potentially observe.

4.2.1 Induced subgraph sampling

Under induced subgraph sampling, we observe edges between the sampled nodes only. Consequently, in our sample we observe $\sum_{a \in S_A} \sum_{b \in S_B} 1_{\{a, b\} \in E}$ edges between distinct categories A and B , out of the maximal number $|S_A| \cdot |S_B|$ we could possibly observe, leading to the trivial estimator

$$\hat{w}(A, B) = \frac{\sum_{a \in S_A} \sum_{b \in S_B} 1_{\{a, b\} \in E}}{|S_A| \cdot |S_B|}. \quad (8)$$

(Note that when S contains the same node multiple times, we count any corresponding sampled edges multiple times as well.)

4.2.2 Star sampling

Under star sampling, on sampling node $a \in A$ we observe the set $E_{a, B} \subset E$ of all edges between a and category $B \neq A$. So we observe $|E_{a, B}|$ edges out of a potential $|B|$ edges between a and B . If we consider all nodes S_A we sampled from A , we observe $\sum_{a \in S_A} |E_{a, B}|$ out of a potential $|S_A| \cdot |B|$ edges. The same applies to nodes S_B sampled in B and their neighbors in A . Consequently, we can estimate the category graph edge weight $w(A, B)$ by dividing the total number of edges we observed between A and B by our estimate of the

maximal number we could potentially observe, *i.e.*,

$$\widehat{w}(A, B) = \frac{\sum_{a \in S_A} |E_{a,B}| + \sum_{b \in S_B} |E_{b,A}|}{|S_A| \cdot |\widehat{B}| + |S_B| \cdot |\widehat{A}|}. \quad (9)$$

Note that because we usually do not know the real sizes of A and B , Eq.(9) uses their estimators $|\widehat{A}|$ and $|\widehat{B}|$. We can employ either Eq.(4) or Eq.(5), as needed.

Observe that the star sampling estimator is potentially more efficient than the trivial induced subgraph estimator, because we include edges (and non-edges) between sampled members of A and B and members of the respective sets that were *not* themselves sampled. For categories with large mean degree, this may represent a substantial increase in information versus the induced subgraph case.

4.3 Population size (N)

In our estimation of category sizes, the population size $N=|V|$ is required. In some cases N is known (*e.g.*, in an OSN context, it may be published by the service provider), but in general this is not the case. Fortunately, where N is not available, we can turn to estimation. For instance, [33] proposes an approach based on a “reversed coupon collector” problem, which can be used with both uniform and non-uniform sampling.

Finally, we note that N is only necessary where absolute values of category sizes are required. Specifically, all edge weights and category sizes can be estimated up to a constant of proportionality without knowing the size of the total population. Thus, if we are interested in ratios of category sizes and/or edge weights (*e.g.*, the relative weight of the A, B connection versus the A, C connection in G^C), then N can be ignored (and replaced by an arbitrary constant in the above equations).

5. NON-UNIFORM SAMPLING

The estimators derived in Section 4 hold under UIS, where every node $v \in V$ is sampled with the same probability. Such a sampling design is rarely feasible in practice. Moreover, in some cases UIS may be also undesirable, *e.g.*, when some categories are irrelevant to our measurement [35].

A more common scenario is *non-uniform* probability sampling, where every node $v \in V$ is sampled with probability proportional to a known weight $w(v)$. Indeed, this is the case for WIS, RW, S-WRW and other principled walk-based sampling methods, provided that samples have adequately converged [20]. Non-uniform samples are by definition biased towards nodes of higher weight (typically degree), which may dramatically distort the estimation results if used without correcting for sampling probabilities [21].

Fortunately, where sampling weights are known (as in the above designs), they can be corrected for by an ap-

propriate (though not necessarily obvious) re-weighting of the measured values. In this section, we rewrite all estimators from Section 4 in such a corrected form.

5.1 Correcting for sample bias

A weighted sample can be unbiased using the Hansen-Hurwitz estimator [25] as shown *e.g.*, in [56,65] for random walks and also used in [51]. Let every node $v \in V$ carry a value $x(v)$. We can estimate the population total $x_{\text{tot}} = \sum_v x(v)$ by

$$\hat{x}_{\text{tot}} = \frac{1}{n} \sum_{v \in S} \frac{x(v)}{\pi(v)}, \quad (10)$$

where $\pi(v)$ is the sampling probability of node v .

In practice, we usually know $\pi(v)$, and thus \hat{x}_{tot} , only up to a constant, *i.e.*, we know the (non-normalized) weights $w(v)$, $w(v) \sim \pi(v)$. Fortunately, we can often address this problem by estimating the ratio of two totals, which makes the unknown constants cancel out. We will use this approach below.

5.2 Estimating category size ($|A|$)

5.2.1 Induced subgraph sampling

Following Eq.(10), we can estimate $|S_A|$ by setting $x(v) \equiv 1_{\{v \in A\}}$. This yields $|\widehat{S}_A| = \frac{1}{n} \sum_{v \in S} \frac{1_{\{v \in A\}}}{\pi(v)} = \frac{1}{n} \sum_{v \in S_A} \frac{1}{\pi(v)}$. Analogously, $|\widehat{S}| = \frac{1}{n} \sum_{v \in S} \frac{1}{\pi(v)}$. Consequently, we can rewrite Eq.(4) as

$$\begin{aligned} |\widehat{A}| &= N \cdot \frac{\sum_{v \in S_A} \frac{1}{\pi(v)}}{\sum_{v \in S} \frac{1}{\pi(v)}} = N \cdot \frac{\sum_{v \in S_A} \frac{1}{w(v)}}{\sum_{v \in S} \frac{1}{w(v)}} \\ &= N \cdot \frac{w_{\cdot 1}(S_A)}{w_{\cdot 1}(S)}, \end{aligned} \quad (11)$$

where

$$w_{\cdot 1}(X) = \sum_{v \in X} \frac{1}{w(v)}$$

is a ‘re-weighted size’ of multiset $X \subset V$.

5.2.2 Star sampling

As in Section 4.1.2, we estimate the size of a category A using Eq.(5), *i.e.*,

$$|\widehat{A}| = N \cdot \widehat{f}_A^{\text{vol}} \cdot \frac{\widehat{k}_V}{\widehat{k}_A}. \quad (12)$$

However, now, the terms $\widehat{f}_A^{\text{vol}}$, \widehat{k}_V and \widehat{k}_A must be calculated taking into account the sampling weights. Indeed, the weighted version of $\widehat{f}_A^{\text{vol}}$ is (after [35])

$$\widehat{f}_A^{\text{vol}} = \frac{1}{\sum_{s \in S} \frac{\deg(s)}{w(s)}} \cdot \sum_{s \in S} \left(\frac{1}{w(s)} \sum_{v \in \mathcal{N}(s)} 1_{\{v \in A\}} \right). \quad (13)$$

Similarly, the estimators Eq.(6) of k_V and k_A can be rewritten respectively by

$$\hat{k}_V = \frac{\sum_{v \in S} \frac{\deg(v)}{w(v)}}{w_{-1}(S)} \quad \text{and} \quad \hat{k}_A = \frac{\sum_{v \in S_A} \frac{\deg(v)}{w(v)}}{w_{-1}(S_A)}. \quad (14)$$

5.3 Estimating category edge weights ($w(A, B)$)

5.3.1 Induced subgraph sampling

Note that in the numerator of Eq.(8), we have a sum over node *pairs*, rather than single nodes. In this case, Hansen-Hurwitz estimator divides every component by the product of weights of the two involved nodes [34], which yields

$$\hat{w}(A, B) = \frac{\sum_{a \in S_A} \sum_{b \in S_B} \frac{1_{\{a, b\} \in E}}{w(a) \cdot w(b)}}{w_{-1}(S_A) \cdot w_{-1}(S_B)}. \quad (15)$$

5.3.2 Star sampling

Finally, under nonuniform sampling, Eq.(9) becomes

$$\hat{w}(A, B) = \frac{\sum_{a \in S_A} \frac{|E_{a, B}|}{w(a)} + \sum_{b \in S_B} \frac{|E_{b, A}|}{w(b)}}{w_{-1}(S_A) \cdot |\hat{B}| + w_{-1}(S_B) \cdot |\hat{A}|}. \quad (16)$$

Again, we have two size estimators Eq.(11) and Eq.(12) to choose from to plug into $|\hat{A}|$ and $|\hat{B}|$. We recommend selecting the one with smaller variance for the specific application. This variance can be estimated, *e.g.*, using bootstrapping [9].

5.4 Sampling via crawling

As we argued in Section 3.1.2, in many online networks the only feasible sampling approach is via crawling. Such techniques result in non-uniform sampling probabilities, and, consequently, sampling weights. For example, under RW the sampling weights converge asymptotically to $w(v) = \deg(v)$ [41]. Using these weights in conjunction with the WIS estimators above allows for consistent estimation of coarse-grained topology from random walk samples.

Of course, consecutive samples collected by crawls are in general correlated, which can potentially affect the efficiency of our estimators. One way to deal with that is to take, say, every T -th sample. For T large enough, this *thinning* technique effectively reduces sample correlations, at a cost of discarding a large portion of available information. Thinning is crucial in some applications, *e.g.*, those primarily based on counting repeated nodes, as in [33]. The ergodicity of standard random walk designs, however, guarantees convergence to the target (WIS) distribution with any effect of autocorrelation vanishing in the limit of sample size. (See Appendix.)

6. SIMULATION RESULTS

6.1 Objective and performance metrics

In this section, we apply our methodology to fully observed graphs from both synthetic and empirical sources. Our objective is to evaluate estimator performance by comparing the (known) values of the category sizes and edge weights in each case with the values inferred using our methods. We use the Normalized Root Mean Square Error (NRMSE) to assess estimation error:

$$\text{NRMSE}(\hat{x}) = \frac{\sqrt{\mathbb{E}[(\hat{x} - x)^2]}}{x}, \quad (17)$$

where x is the real value and \hat{x} is the estimate.

6.2 Generated topologies

First, we consider synthetic graphs. By simulating G , we control many crucial parameters (such as graph density, or category size and tightness) and study the effect of these parameters on the efficiency of our estimators.

6.2.1 Graph model

We consider a graph G with $N = 88,850$ nodes partitioned into 10 categories. Their sizes range from $|C|=50$ to $|C|=50000$. Initially, nodes in each category form a k -regular random graph, with the average node degree ranging from $k=5$ to $k=49$. In addition, we add $N \cdot k/10$ random edges between nodes in different categories. The resulting graph G is connected (in all instances we used) and has $|E| = 0.6 \cdot N \cdot k$ edges. By construction, G has a very strong community structure. In order to study the effect of community tightness, we next permute randomly the category labels of a fraction $\alpha \in [0, 1]$ of nodes. For $\alpha=0$, node categories follow the strong community structure, whereas for $\alpha=1$ the categories are completely independent of the graph structure.

6.2.2 Category sizes

We first study the efficiency of the category size estimators, Eq.(4) and Eq.(5). We present the results in the top row of Fig. 3 and make the following observations.

In all of our simulated cases, all estimators converge to the true value as sample size increases. Moreover, the star estimator performs better than the induced subgraph estimator, although its efficiency can depend on properties of G . For example, (i) the denser the graph, the better the star estimator is (Fig. 3(a)), but (ii) its efficiency can be limited when clustering closely follows the category structure (Fig. 3(b)). In contrast, the induced subgraph estimator is not affected by any of these properties. We also observe that both estimators perform better for larger categories (Fig. 3(c)). In Fig. 3(d), we show the CDF of the NMSE of all (ten) estimators of the category sizes.

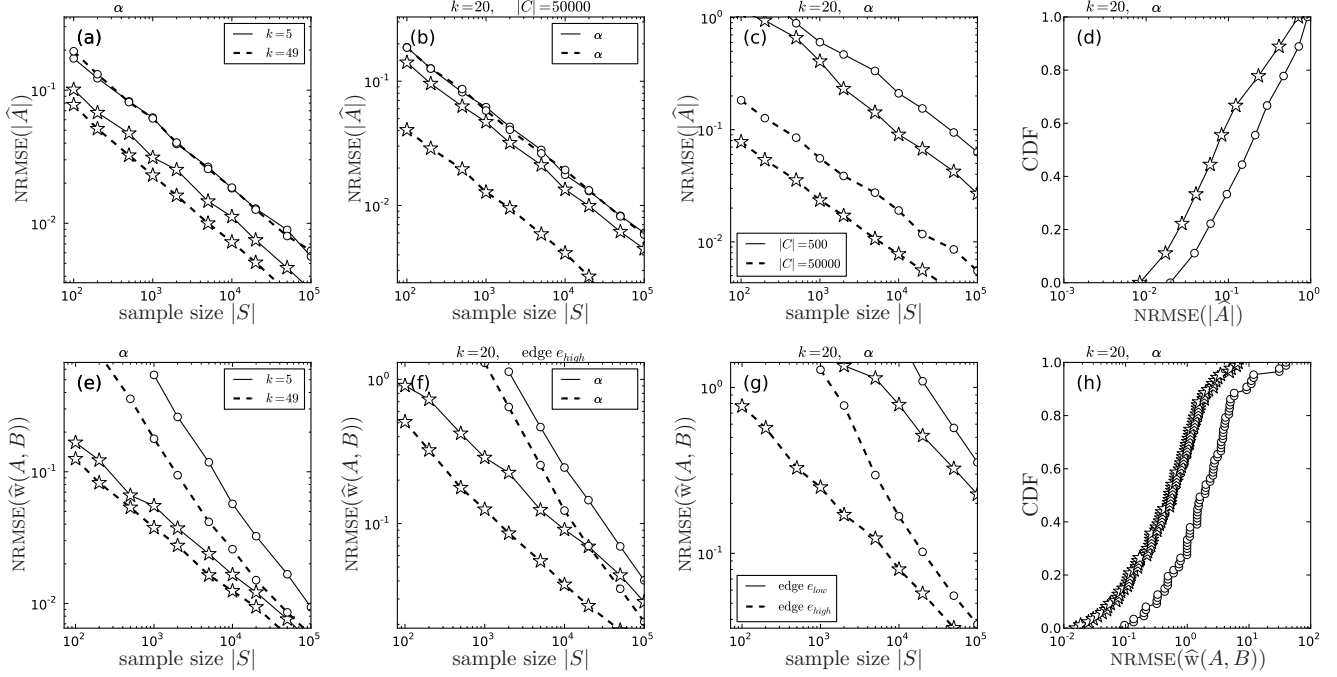


Figure 3: Simulations of UIS on synthetic graphs. We estimate category sizes (top) and category edge weights (bottom), using induced subgraph sampling (circles) and star sampling (stars).

Dataset	$ V $	$ E $	k_V
Facebook: Texas [62]	36 364	1 590 651	87.5
Facebook: New Orleans [64]	63 392	816 885	25.8
P2P [40]	62 561	147 877	4.7
Epinions [54]	75 877	405 738	10.7

Table 1: Empirical topologies used in Sec. 6.3.

6.2.3 Category edge weights

In the bottom row of Fig. 3, we use Eq.(8) and Eq.(9) to estimate the category edge weights under induced and star sampling designs, respectively.

Again, both estimators converge, with the star estimator performing better than the induced one. As before, the star estimator benefits from higher graph density (Fig. 3(e)) and looser category structure (Fig. 3(f)). However, in this case the induced estimator is affected by these properties as well. Finally, in Fig. 3(g) we compare the estimation efficiency of low-weight edge e_{low} (defined as the edge with 25th percentile weight) with the high-weighted edge e_{high} (75th percentile). As before, both estimators perform better for large estimated values.

6.3 Empirically observed topologies

6.3.1 Datasets

We consider four fully known topologies described in Table 1. We use two graphs extracted from Facebook because (i) they significantly differ in density, and

(ii) Facebook is our focus in the experimental study of Section 7.

In Section 6.2, we have seen that star sampling performs the worst if categories are aligned with the communities (dense clusters) existing in graphs. We decided to simulate presumably the worst-case category partition from the star sampling point of view. In particular, we use a standard community finding algorithm based on eigenvalues [47] to identify the 50 largest communities, and define each such community to be a category. All the remaining smaller categories (if any) are then grouped together as the 51st category.

From these known graphs we then generate synthetic datasets by three different sampling methods: UIS, RW and S-WRW. Under S-WRW [35], we use equal category weights for all categories, and we set $\tilde{f}_\Theta = 0$ (because there are no irrelevant categories) and $\gamma = \infty$ (for simplicity). As previously, our interest is in whether our estimators (applied to these realistic samples) will accurately reconstruct the true properties of the graphs in question.

6.3.2 Category size

We study the efficiency of the category size estimators in the top row of Fig. 4. Due to lack of space, we only report the median NRMSE across all categories. In Fig. 3(d), this would correspond to the points on the horizontal line $Y = 0.5$.

The main observation is that, in contrast to Sec-

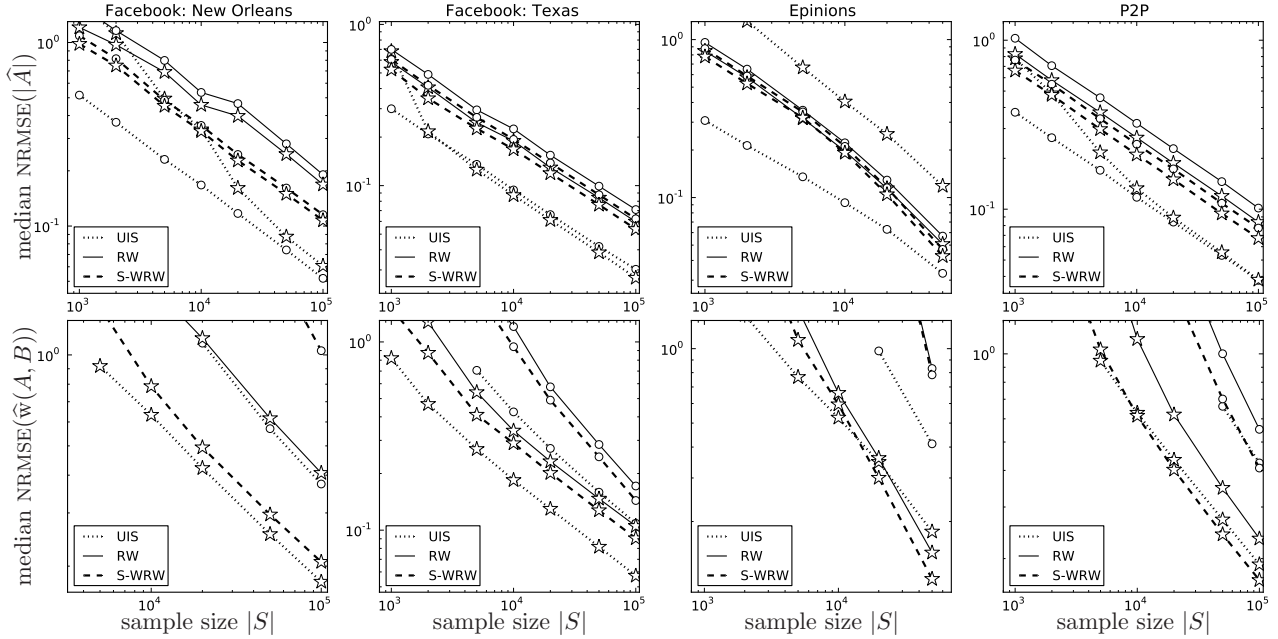


Figure 4: Simulations on empirically observed graphs. We estimate category sizes (top) and category edge weights (bottom), using induced subgraph sampling (circles) and star sampling (stars).

tion 6.2, here the induced estimators can outperform the star estimators. This is particularly visible under UIS, probably because of the highly skewed node degree distributions. Such a distribution increases the variance of the average degree estimator \hat{k}_A that is used in the star-based size estimation in Eq.(5).⁴

However, in contrast to UIS, under RW and S-WRW star sampling usually performs better. This can also be explained by the highly skewed node degree distribution. Indeed, because both RW and S-WRW visit high-degree nodes more often than UIS, their star samples inherently collect and exploit more information about neighbor categories, which translates to a better performance. This effect is similar to the better star sampling performance under higher graph density in Section 6.2.

6.3.3 Category edge weights

While there is no clear winner in the category size estimation, in the category edge weight estimation star sampling consistently and significantly outperforms induced sampling. Indeed, in Fig. 4(e-h), the induced estimators often need 5-10 times more samples to achieve

⁴We might address this problem by modifying Eq.(5) to take *e.g.*, $\hat{k}_A = \hat{k}_V$ or a similar model-based extension. Such modifications may greatly reduce the variance of size estimation, albeit at the cost of some bias. (Indeed, this is an example of the classic “precision vs accuracy” tradeoff.) Note that such modifications can allow us to use Eq.(5) to estimate $|A|$, even if none of our sampled vertices were drawn from A . Our initial experiments with such modifications have been encouraging, but we do not treat them in depth here.

Dataset	Studied categories	Crawl type	% cat. samples	# total samples
2009 [20]	Regional (507) (34% of population)	MHRW09	34%	28x81K
		RW09	41%	28x81K
		UIS09	34%	28x35K
2010 [35]	Colleges (10K+) (3.5% of population)	RW10	9%	25x40K
		S-WRW10	86%	25x40K

Table 2: Facebook datasets.

the same accuracy as star estimators.

UIS clearly performs best, especially when estimating category sizes. Not surprisingly, direct independence sampling should be preferred whenever available. In the more practical scenarios, however, we are limited to exploration-based techniques. In our simulations, S-WRW is consistently better than RW. Note that because all categories (and thus nodes) are relevant, this advantage of S-WRW is purely due to stratification. Moreover, the advantage of S-WRW increases with higher heterogeneity of category sizes (not shown here), which is in agreement with [35].

7. FACEBOOK CATEGORY GRAPHS

In this section, we use the estimators developed in this paper to infer several category graphs from Facebook.

7.1 Data sets

In our previous work [20,35], we collected samples of Facebook users (about 10.1 million total users), with publicly available information. These datasets are summarized in Table 2 and are used as input for the esti-

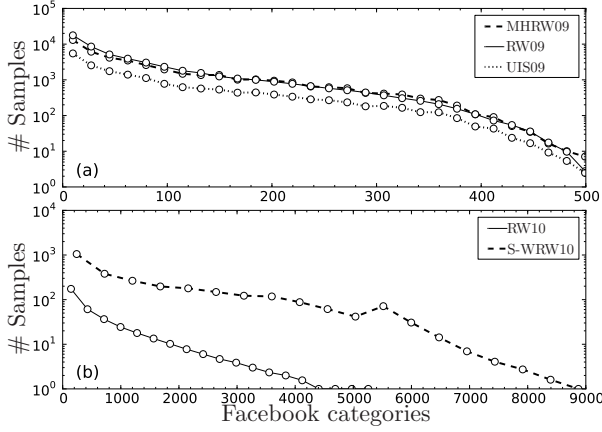


Figure 5: Number of samples per category in Facebook 2009 (top) and 2010 (bottom).

mators of this paper. These datasets were collected using HTML scraping, which allowed us to collect for each user v not only v 's category, but also the list of v 's friends together with their categories; *i.e.*, we effectively collected a star sample of Facebook users. By discarding the information about v 's nodes, we can also use the induced subgraph estimators, for comparison.

The 2009 data sets: These data sets were collected in April 2009 [20], using three existing sampling techniques, UIS, MHRW and RW, as summarized in Table 2. At that time, a Facebook user could be a member of any of four different types of categories, called “networks” in the Facebook terminology. Three of them, *high school*, *college* and *workplace*, required passing a verification process, usually based on an email account from the institution in question. The fourth category, *geographical region*, did not require any verification, and indicated the user's city, state or country. In this paper, we consider the geographical region categories from the 2009 data sets. Each dataset consisted of 100-1000 samples from each of the 507 geographical regions, as shown in Fig. 5(a); UIS collected about two times fewer samples than the other two techniques.

The 2010 data sets: The geographical region category was phased out in June 2009. Therefore, the data sets we collected in 2010 [35] contain only the three remaining categories, from which we chose colleges as the category studied in this paper. Furthermore, Facebook switched from 32 bit to 64 bit userIDs, thus leading to a sparse userID space, which made UIS impractical to apply. For this reason, in our 2010 Facebook data sets we collected only a RW sample (because RW proved to outperform MHRW [20,51]) as well as three variants of S-WRW [35]. A full length (1M) RW typically collected only 0-10 samples of a particular college (Fig. 5(b)). This is because of a relatively small college population (about 3.5%) and a large number of colleges (more than

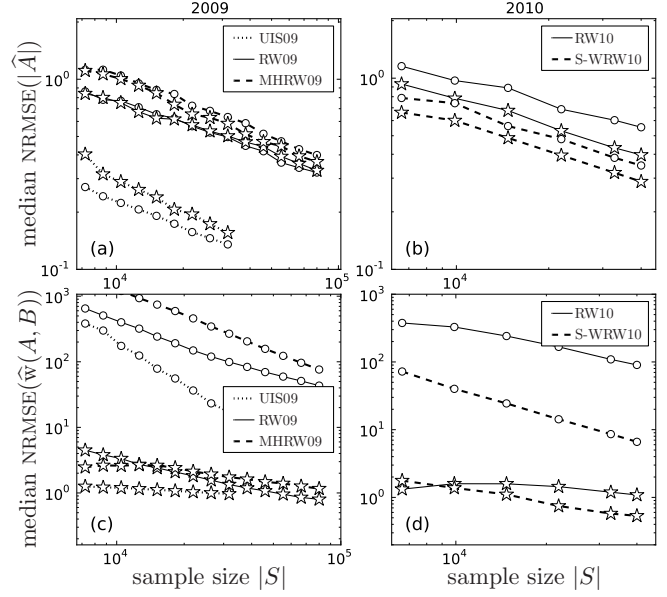


Figure 6: Results for 100 most popular regional networks in 2009 (a,c) and 100 college networks in 2010 (b,d): category size estimation (a,b), and edge weight estimation (c,d).

10,000). Fortunately, S-WRW, a technique designed to oversample particular categories (here colleges), improves that result by at least one order of magnitude.

7.2 Category graph estimation

We present our results in Fig. 6. To calculate NRMSE we use as ground truth the average of estimation over all samples for each crawl type. In addition, we treat each of the 28 and 25 different walks, for the 2009 and 2010 data sets respectively, as a different sample.

7.2.1 Category size

We show the results of Facebook category size estimation in Fig. 6(a,b). Similarly to what we observed in the simulations in Section 6, UIS performs the best, and S-WRW outperforms RW. MHRW performs the worst, which was also expected given the recent studies of MHRW in [20,51]. Under UIS, the induced estimator performs better. Under RW and S-WRW, the star version is better, especially when categories are small, as in the 2010 data set.

7.2.2 Category edge weight

The estimation of category edge weights in Facebook, shown in Fig. 6(c,d), also confirms the observations in the simulations of Section 6. Indeed, all star estimators dramatically outperform their induced counterparts. And, as before, the sampling techniques ordered from the best to worst are: UIS, S-WRW, RW and MHRW.

Finally, note that NRMSEs in Fig. 6(a-d) are relatively high, even under star (*i.e.*, the better performing) sampling. This is because these plots reach only relatively small sample sizes $|S|$ (*i.e.*, 25 or 28 times smaller than the entire sample at our disposal). Therefore, one could extrapolate the plots in Fig. 6 by much more than a decade to the right, further reducing the values of NRMSE. Moreover, in the data sets that we eventually prepare, we combine together several outcomes of different, independent sampling techniques, which should further limit the estimation variance. Therefore, the results in Fig. 6 should be treated as a guideline about the relative efficiency of the sampling techniques, rather than a comparison of the the absolute values of NRMSE.

7.3 Geosocial visualization

Finally, we have developed a highly customizable, web-based tool for visualization of our Facebook category graphs. We have made a beta-version of the tool available at www.geosocialmap.com and invite the reader to use it to experiment with the category-graphs described in this paper. This can be used to gain insight into the friendship relations among these categories, as defined *in Facebook*.⁵

7.3.1 Cross-country friendships

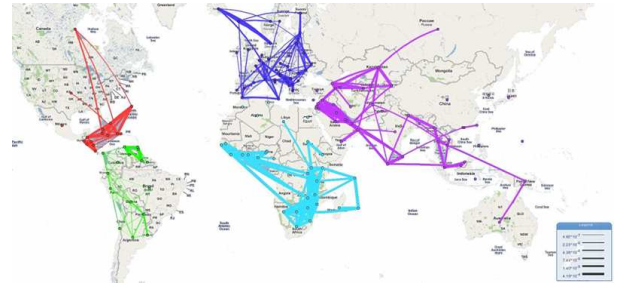
As mentioned earlier in Section 7.1, the 2009 data set contains the geographical region information, at various granularities depending on Facebook’s penetration in that region. This may be either a user’s city or state, (*e.g.*, for USA, Canada, UK) or the entire country (more typically).

As an example, we create the country-to-country friendship graph. To this end, we first merged together all categories coming from the same country. Next, we estimated the sizes of the resulting categories. Because, according to Fig. 6(a), the UIS induced sampling performed exceptionally well, we used it in the category size estimation. This information was next fed to the star estimators of category edge weights. Finally, for every edge, we take the average of the three estimates (resulting from UIS, MHRW and RW). Fig. 7(a) presents a subset of “The world according to Facebook” graph.

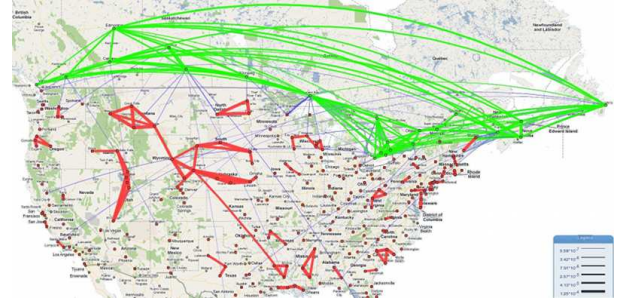
7.3.2 North America

For the USA and Canada, the 2009 data set contains the geographical information at the granularity of 272 counties and provinces. This allows us to create the

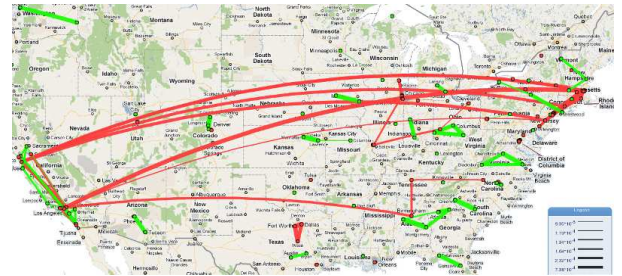
⁵However, one should be careful about declaring categories in Facebook as representative of the real world. First, Facebook attracts some age groups more than others. Second, many Facebook users do not declare (or hide) their category membership. Finally, a user might have mistakenly chosen her category. For example, the third strongest link for “Greece” is “Athens, GA”, which is clearly mistaken for Athens, Greece.



(a) Intra-continental country connections: Note the strong cliques formed between Middle Eastern countries and South-East-Asian countries. There is no Facebook in China.



(b) North-American regions: Physical distance is a major factor in the United States (red), but seemingly less so in Canada (green). Additionally, US and Canada are relatively weakly interconnected (thin blue lines).



(c) Top 133 US colleges according to the “US News World Report’09”: Physical distance is a major factor for public colleges (green), but seemingly less so for private ones (red).

Figure 7: The friendship graph between regional networks. Available at www.geosocialmap.com

North American friendship map. We followed the same steps as in Section 7.3.1. An example is presented in Fig. 7(b).

7.3.3 US colleges

Both the 2009 and 2010 data sets contain college categories. We chose the 2010 data set to create a college-to-college friendship graph. This data set consists of one RW sample and three S-WRW ones. Because S-WRW performed much better than RW (see Fig. 6(b,d)), we decided to use the three S-WRW samples only. Moreover, this time we estimated the size with the help

of the star estimators, because they performed better (Fig. 6(b)). Finally, as before, we fed the resulting category sizes into the star estimators of category edge weights, and we averaged the three S-WRW estimates into a final estimate. Fig. 7(c) presents a subgraph of the resulting category graph.

8. RELATED WORK

Node sampling in graphs. Most state-of-the-art crawling-based node sampling techniques use variants of random walks (RW), such as the classic RW [20,27,41,51,56], Metropolis-Hasting RW (MHRW) [18,20,42,51,60], multiple dependent RW [52], multigraph RW [19], RW with jumps [6,30,38,53], and weighted RW [35]. Based on the resulting (uniform or non-uniform) sample of nodes, there exist principled methods to estimate local graph properties (degree distribution, assortativity and clustering coefficient). [34] is an excellent introduction; other examples include [3,6,20,21,26,37,51–53,59]. In our prior work [19,20,35], we used random-walk based crawls to collect user samples, which we use as input to the estimators proposed in this paper.

Topology inference. Much classic work on inference for basic network properties from node samples was done by Ove Frank and colleagues; see particularly [13–15], which introduce Horvitz-Thompson estimators of edge totals (*i.e.*, volumes) from probability samples of nodes. Early results involving topology inference from induced subgraph and star sampling were reviewed by [16]. This prior work focused on the case of known population and category sizes, and assumed without-replacement designs.

Breadth First Search (BFS) has been used to sample topology *e.g.*, in [4,43,44]. However, a BFS sample is known to introduce a strong bias towards high degree nodes [7,20,36,37,46,70], which makes it not representative with respect to many metrics. Although this degree bias can often be significantly corrected for [36], the BFS sample covers only the neighborhood of the arbitrary starting node, which is not necessarily representative of the entire topology.

[38] evaluates a number of sampling methods and the graphs they induce. The authors conclude that Forest Fire [39], intuitively a hybrid of RW and BFS, produces topology samples that resemble the original graph the most. However, Forest Fire is subject to the same biases as BFS described above.

Another approach for inferring network structure is matrix completion of the distance matrix [10,68]. However, this approach faces its own challenges when applied to OSN samples. First, the distance matrix is typically high rank and one has to carefully identify a low rank structure [10]. Second, unlike traceroutes or tomographic techniques, crawling does not yield a random sample of distances [10,68].

Induced subgraph vs. star sampling [34] is a good summary of these two sampling designs. Induced subgraph sampling has been studied, *e.g.*, in [34,37,38] Star sampling is similar to egonet sampling [66], except that under star sampling we do not see edges between neighbors of a sampled node. Our contribution here is to apply these measurement schemes in the context of category graph estimation.

Block models and mixing rates The use of partitions to produce reduced-form versions of larger networks has an extensive history in the social network literature, primarily under the label of “block modeling;” see [49,66] for extensive reviews. Block models with known partitions are sometimes called “confirmatory” block models, and have been studied largely from a statistical point of view (*e.g.*, [11,12] and [66] ch. 16). Much of the latter interest is in modeling the edge weights (“block densities” or “mixing rates”) from covariates or other information in a fully-observed context, with considerable additional interest in the case where the network is observed but the categories are latent [48,58]. Estimation of mixing rates from uniform node samples for categories of known size is also a well-known problem (*see.*, *e.g.*, [13,45,69]). Comparable methods for link-trace samples are less well-developed, though see [17,24,27,28].

Although estimation of mixing rates from sampled data is relatively straightforward where categories are of known size and the number of categories $|\mathcal{C}|$ is fairly small (so that a random sample provides large numbers of vertex pairs in each pair of categories), it is much more difficult when $|\mathcal{C}|$ is large and category sizes are not known. Our techniques thus extend the prior literature on block models and mixing rates to cases such as group interaction in OSNs and other large-scale social networks, in which one must estimate interaction among many groups of uncertain size from (typically non-uniformly) sampled data. Our work also differs from much recent social network literature in being design-based rather than model-based; design-based inference is frequently easier to employ than model-based inference, although both approaches have merits [61].

Facebook colleges. The Facebook social graph has been measured and studied in the past. For example, [22] studies the interactions between all 4.2M Facebook users in 492 universities in North America between Feb 2004 and March 2006. (As a side note, the interpretation is hindered by the full anonymization of user and universities.) [62] studies the social structure within 100 Facebook college categories. Given the above full datasets, one could apply Eq.(3) and create the category graph. In contrast, our methodological contribution lies in estimating the category graph from a sample of nodes, not from the fully known user graph.

Social graph visualization. There exist many tools that visualize social graphs (including Facebook), for example [1,2,29]. www.geosocialmap.com differs from most of these tools in that it (i) is category-centric (vs user-centric), (ii) contains an aggregated information view of *entire* Facebook population, (iii) is well suited for data exploration (*e.g.*, allows arbitrary selection of categories), and (iv) accepts as input any weighted graph with arbitrary set of node/edge attributes (ongoing work).

9. CONCLUSION

Estimation performance. In this paper, we derive a number of category graph estimators for probability samples of nodes, uniform (Section 4) and non-uniform (Section 5). We evaluate their performance in simulation (Section 6) and on Facebook samples (Section 7). We showed that they all converge to their true values for reasonable sample sizes, a result we extend formally in the Appendix. Based on our evaluation, we also provide recommendations, summarized as follows. When estimating category sizes, there is no universal choice between induced and star sampling. For example, the performance of the star estimator improves (i) in dense graphs, (ii) in graphs with homogeneous node degree distribution, (iii) in graphs with weaker community structure, and (iv) under sampling techniques that oversample high degree nodes. In contrast, when estimating the category edge weights, the star estimators are a clear winner; the induced subgraph estimators often need 5-10 times more samples to achieve the same accuracy. Finally, the sampling techniques strongly affect estimator efficiency. They can be ordered from best to worst as follows: UIS, S-WRW, RW and MHRW.

Potential applications. We applied our methodology to samples of Facebook users and we estimated potentially interesting category graphs, such as the global friendship map, or the friendship network of US colleges. We visualized and made publicly available these weighted topologies at www.geosocialmap.com.

In addition to purely descriptive uses, the techniques described here can also be employed as a first step towards model-based analysis. Using the unnormalized edge weights together with the number of possible edges within each cut yields the numbers of possible and realized edges needed for likelihood-based analysis of interaction probabilities. For instance, given additional features associated with each category (*e.g.*, for universities, their size, location, ranking, and expense), one can model the inter-category mixing rates as a function of category features (*e.g.*, the effect of geographical distance on tie probability). This permits both hypothesis testing for putative theories of tie formation and ex ante prediction of interaction rates among new

or unobserved categories (given their hypothesized features) for extremely large, incompletely observed networks. Given the large and growing literature on statistical modeling of networks (*e.g.*, [8,23,31,32,50,55,63] among many others), the potential for applications in this area is substantial.

10. REFERENCES

- [1] My Friend Map (<http://www.facebook.com/myfriendmap>).
- [2] Touchgraph (<http://www.touchgraph.com/facebook>).
- [3] N. Ahmed, J. Neville, and R. Kompella. Reconsidering the Foundations of Network Sampling. In *WIN*, 2010.
- [4] Y. Ahn, S. Han, H. Kwak, S. Moon, and H. Jeong. Analysis of topological characteristics of huge online social networking services. In *WWW*, pages 835–844, 2007.
- [5] D. Aldous and J. A. Fill. *Reversible Markov Chains and Random Walks on Graphs*. In preparation.
- [6] K. Avrachenkov, B. Ribeiro, and D. Towsley. Improving Random Walk Estimation Accuracy with Uniform Restarts. In *17th Workshop on Algorithms and Models for the Web Graph*, 2010.
- [7] L. Becchetti, C. Castillo, D. Donato, and A. Fazzzone. A comparison of sampling techniques for web graph characterization. In *LinkKDD*, 2006.
- [8] C. T. Butts. Permutation models for relational data. *Sociological Methodology*, 37(1):257–281, 2007.
- [9] B. Efron and R. Tibshirani. *An Introduction to the Bootstrap*. Chapman and Hall, London, 1993.
- [10] B. Eriksson, P. Barford, J. Sommers, and R. Nowak. DomainImpute: Inferring Unseen Components in the Internet. In *IEEE INFOCOM Mini-Conference*, 2011.
- [11] S. E. Fienberg, M. M. Meyer, and S. Wasserman. Statistical analysis of multiple sociometric relations. *Journal of the American Statistical Association*, 80:51–67, 1985.
- [12] S. E. Fienberg and S. Wasserman. Categorical data analysis of single sociometric relations. *Sociological Methodology*, 12:156–192, 1981.
- [13] O. Frank. Estimation of graph totals. *Scandinavian Journal of Statistics*, 4(2):81–89, 1977.
- [14] O. Frank. Survey sampling in graphs. *Journal of Statistical Planning and Inference*, 1:235–264, 1977.
- [15] O. Frank. Sampling and estimation in large social networks. *Social Networks*, 1(1):91–101, 1978.
- [16] O. Frank. Random sampling and social networks: A survey of various approaches. *Mathematiques, Informatique, et Sciences Humaines*, 26:19–33, 1988.
- [17] K. J. Gile and M. S. Handcock. Respondent-driven sampling: An assessment of current methodology. *Sociological Methodology*, 40, 2010.
- [18] W. R. Gilks, S. Richardson, and D. J. Spiegelhalter. *Markov Chain Monte Carlo in Practice*. Chapman and Hall/CRC, 1996.
- [19] M. Gjoka, C. T. Butts, M. Kurant, and A. Markopoulou. Multigraph Sampling of Online Social Networks. *To appear in JSAC on Measurement of Internet Topologies*, 2011.
- [20] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou. Walking in Facebook: A Case Study of Unbiased Sampling of OSNs. In *INFOCOM*, 2010.
- [21] M. Gjoka, M. Kurant, C. T. Butts, and A. Markopoulou. Practical Recommendations on Sampling OSN Users by Crawling the Social Graph. *To appear in JSAC on Measurement of Internet Topologies*, 2011.
- [22] S. A. Golder, D. Wilkinson, and B. A. Huberman. Rhythms of Social Interaction: Messaging within a Massive Online Network. In *3rd International Conference on Communities and Technologies*, 2007.
- [23] S. M. Goodreau, J. A. Kitts, and M. Morris. Birds of a feather, or friend of a friend?: Using exponential random graph models to investigate adolescent social networks. *Demography*, 46(1):103–125, 2009.

- [24] M. S. Handcock and K. J. Gile. Modeling networks with sampled data. *Annals of Applied Statistics*, 4(1):5–25, 2010.
- [25] M. Hansen and W. Hurwitz. On the Theory of Sampling from Finite Populations. *Annals of Mathematical Statistics*, 14(3), 1943.
- [26] S. J. Hardiman, P. Richmond, and S. Hutzler. Calculating statistics of complex networks through random walks with an application to the on-line social network Bebo. *The European Physical Journal B*, 71(4):611–622, Aug. 2009.
- [27] D. D. Heckathorn. Respondent-Driven Sampling: A New Approach to the Study of Hidden Populations. *Social Problems*, 44:174–199, 1997.
- [28] D. D. Heckathorn. Respondent-driven sampling ii: Deriving valid population estimates from chain-referral samples of hidden populations. *Social Problems*, 49(1):11–34, 2002.
- [29] J. Heer and D. Boyd. Vizster: Visualizing online social networks. *IEEE InfoVis*, 2005.
- [30] M. R. Henzinger, A. Heydon, M. Mitzenmacher, and M. Najork. On near-uniform URL sampling. In *WWW*, 2000.
- [31] P. D. Hoff, A. E. Raftery, and M. S. Handcock. Latent space approaches to social network analysis. *Journal of the American Statistical Association*, 97(460):1090–1098, 2002.
- [32] D. R. Hunter and M. S. Handcock. Inference in curved exponential family models for networks. *Journal of Computational and Graphical Statistics*, 15:565–583, 2006.
- [33] L. Katzir, E. Liberty, and O. Somekh. Estimating Sizes of Social Networks via Biased Sampling. In *WWW*, 2011.
- [34] E. D. Kolaczyk. *Statistical Analysis of Network Data*, volume 69 of *Springer Series in Statistics*. Springer New York, 2009.
- [35] M. Kurant, M. Gjoka, C. T. Butts, and A. Markopoulou. Walking on a Graph with a Magnifying Glass: Stratified Sampling via Weighted Random Walks. In *Sigmetrics*, 2011.
- [36] M. Kurant, A. Markopoulou, and P. Thiran. Towards Unbiased BFS Sampling. *To appear in JSAC on Measurement of Internet Topologies*, (December):1–14, 2011.
- [37] S. H. Lee, P.-J. Kim, and H. Jeong. Statistical properties of Sampled Networks. *Phys. Rev. E*, 73:16102, 2006.
- [38] J. Leskovec and C. Faloutsos. Sampling from large graphs. In *KDD*, pages 631–636, 2006.
- [39] J. Leskovec, J. Kleinberg, and C. Faloutsos. Graphs over time: densification laws, shrinking diameters and possible explanations. In *KDD*, 2005.
- [40] J. Leskovec, J. Kleinberg, and C. Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):2, Mar. 2007.
- [41] L. Lovász. Random walks on graphs: A survey. *Combinatorics, Paul Erdős is Eighty*, 2(1):1–46, 1993.
- [42] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equation of state calculation by fast computing machines. *Journal of Chemical Physics*, 21:1087–1092, 1953.
- [43] A. Mislove, H. S. Koppula, K. P. Gummadi, P. Druschel, and B. Bhattacharjee. Growth of the Flickr social network. In *WOSN*, 2008.
- [44] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhattacharjee. Measurement and analysis of online social networks. In *IMC*, pages 29–42, 2007.
- [45] M. Morris. A log-linear modeling framework for selective mixing. *Mathematical Biosciences*, 107:349–377, 1991.
- [46] M. Najork and J. L. Wiener. Breadth-first search crawling yields high-quality pages. In *WWW*, 2001.
- [47] M. Newman. Finding community structure in networks using the eigenvectors of matrices. *Physical Review E*, 74(3):1–19, Sept. 2006.
- [48] K. Nowicki and T. A. B. Snijders. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96:1077–1087, 2001.
- [49] M. Nunkesser and D. Sawitzki. In U. Brandes and T. Erlebach, editors, *Network Analysis: Methodological Foundations*, chapter 10, pages 253–291. Springer-Verlag, Berlin, 2005.
- [50] P. E. Pattison, S. Wasserman, G. L. Robins, and A. M. Kanfer. Statistical evaluation of algebraic constraints for social networks. *Journal of Mathematical Psychology*, 44:536–568, 2000.
- [51] A. Rasti, M. Torkjazi, R. Rejaie, N. Duffield, W. Willinger, and D. Stutzbach. Respondent-driven sampling for characterizing unstructured overlays. In *Infocom Mini-conference*, pages 2701–2705, 2009.
- [52] B. Ribeiro and D. Towsley. Estimating and sampling graphs with multidimensional random walks. In *IMC*, 2010.
- [53] B. Ribeiro, P. Wang, and D. Towsley. On Estimating Degree Distributions of Directed Graphs through Sampling. *UMass Technical Report*, 2010.
- [54] M. Richardson, R. Agrawal, and P. Domingos. Trust management for the semantic web. *The SemanticWeb-ISWC 2003*, pages 351–368, 2003.
- [55] G. L. Robins and M. Morris. Advances in exponential random graph (p^*) models. *Social Networks*, 29:169–172, 2007.
- [56] M. Salganik and D. D. Heckathorn. Sampling and estimation in hidden populations using respondent-driven sampling. *Sociological Methodology*, 34(1):193–240, 2004.
- [57] T. A. Severini. *Elements of Distribution Theory*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge, UK, 2005.
- [58] T. A. B. Snijders and K. Nowicki. Estimation and prediction for stochastic block models for graphs with latent block structure. *Journal of Classification*, 14:75–100, 1997.
- [59] M. Stumpf, C. Wiuf, and R. May. Subnets of scale-free networks are not scale-free: sampling properties of networks. *Proceedings of the National Academy of Sciences of the United States of America*, 102(12):4221, Mar. 2005.
- [60] D. Stutzbach, R. Rejaie, N. Duffield, S. Sen, and W. Willinger. On unbiased sampling for unstructured peer-to-peer networks. In *IMC*, 2006.
- [61] S. K. Thompson. *Sampling*. Wiley, New York, second edition, 2002.
- [62] A. Traud, P. Mucha, and M. Porter. Social Structure of Facebook Networks. *Arxiv preprint arXiv:1102.2166*, 2011.
- [63] M. A. J. van Duijn, T. A. B. Snijders, and B. H. Zijlstra. p_2 : A random effects model with covariates for directed graphs. *Statistica Neerlandica*, 58:234–254, 2004.
- [64] B. Viswanath, A. Mislove, M. Cha, and K. Gummadi. On the evolution of user interaction in facebook. In *WOSN*, volume 09, pages 37–42, 2009.
- [65] E. Volz and D. D. Heckathorn. Probability based estimation theory for respondent driven sampling. *Journal of Official Statistics*, 24(1):79–97, 2008.
- [66] S. Wasserman and K. Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, Cambridge, 1994.
- [67] W. Willinger, R. Rejaie, M. Torkjazi, M. Valafar, and M. Maggioni. OSN Research: Time to Face the Real Challenges. In *HotMetrics*, 2009.
- [68] W. Xu, E. Mallada, and A. Tang. Compressive Sensing over Graphs. *INFOCOM*, 2011.
- [69] K. Yamaguchi. Homophily and social distance in the choice of multiple friends: An analysis based on conditionally symmetric log-bilinear association model. *Journal of the American Statistical Association*, 85(410):356–366, 1990.
- [70] S. Ye, J. Lang, and F. Wu. Crawling online social graphs. In *Asia-Pacific Web Conference (APWEB)*, pages 236–242, 2010.

Appendix: Consistency of the estimators

A desirable property of a statistical estimator is that of *consistency*. A statistical estimator (X_n) is a function of the sample size $(n = |S|)$, and is said to be *consistent* if it converges in probability (\xrightarrow{P}) to the true value of interest (θ) [57] (which also implies *asymptotic unbiasedness*). Formally: If $X_n \xrightarrow{P} \theta$, as $n \rightarrow \infty$, then X_n is said to be consistent for θ . To prove the consistency of the estimators in this paper we invoke two classic theorems in probability: (1) The Law of Large Numbers, and (2) Slutsky's Theorem⁶; which require the following assumptions: For the uniform case we need to assume that the mean and variance are finite ($\theta < \infty$; $\sigma^2 < \infty$); for the non-uniform case we need to make an additional assumption on the sampling weights so as to guarantee the consistency of the Hansen-Hurwitz (HH) estimator, specifically that the sum of the weights be bounded ($\sum_{v \in V} w(v) \leq c$)⁷. Both of these conditions are satisfied for finite graphs.

LLN and Slutsky's Theorem

THEOREM 10.1 (Law of Large Numbers (LLN)). Let X_1, X_2, \dots be i.i.d. random variables with $EX_i = \theta$ and $\text{Var } X_i = \sigma^2 < \infty$. Then $\bar{X}_n = \frac{1}{n} \sum_{v \in S} X(v) \xrightarrow{P} \theta$.

THEOREM 10.2 (Slutsky's Theorem). Let $X_n \xrightarrow{P} \alpha$ and $Y_n \xrightarrow{P} \beta$, where α and β , respectively, are real numbers. Then

$$(p.1) \quad X_n + Y_n \xrightarrow{P} \alpha + \beta;$$

$$(p.2) \quad X_n \cdot Y_n \xrightarrow{P} \alpha \cdot \beta;$$

$$(p.3) \quad \frac{X_n}{Y_n} \xrightarrow{P} \frac{\alpha}{\beta}, \text{ where } \beta \neq 0.$$

Uniform sampling estimators

$$\text{Eq.(4): } |\hat{A}| = N \cdot \frac{|S_A|}{|S|} = \frac{1}{n} \sum_{v \in S} 1_{\{v \in A\}} \xrightarrow{P} |A| \text{ by the LLN.}$$

$$\text{Eq.(6): } \hat{k}_V = \frac{\sum_{v \in S} \deg(v)}{|S|} \xrightarrow{P} k_V \text{ and}$$

$$\hat{k}_A = \frac{\sum_{v \in S_A} \deg(v)}{|S_A|} \xrightarrow{P} k_A \text{ by the LLN (as above).}$$

$$\begin{aligned} \text{Eq.(7): } \hat{f}_A^{\text{vol}} &= \frac{1}{\text{vol}(S)} \sum_{s \in S} \sum_{v \in \mathcal{N}(s)} 1_{\{v \in A\}} \\ &= \frac{\frac{1}{n} \sum_{s \in S} \sum_{v \in \mathcal{N}(s)} 1_{\{v \in A\}}}{\frac{1}{n} \text{vol}(S)} \xrightarrow{P} f_A^{\text{vol}} \text{ by an application} \\ &\text{of the LLN to both the numerator and denominator, separately, followed by an application of Slutsky's Theorem (p.3).} \end{aligned}$$

$$\text{Eq.(5): } |\hat{A}| = N \cdot \hat{f}_A^{\text{vol}} \cdot \frac{\hat{k}_V}{\hat{k}_A} \xrightarrow{P} |A| \text{ by two applications of Slutsky's Theorem (p.2 and p.3) and consistency of the individual estimators.}$$

⁶For more details about these two theorems see [57].

⁷There are some alternate assumptions on the weights that can be made to guarantee convergence.

$$\begin{aligned} \text{Eq.(8): } \hat{w}(A, B) &= \frac{\sum_{a \in S_A} \sum_{b \in S_B} 1_{\{a, b\} \in E}}{|S_A| \cdot |S_B|} \\ &= \frac{\frac{N^2}{n^2} \sum_{a \in S_A} \sum_{b \in S_B} 1_{\{a, b\} \in E}}{\frac{N^2}{n^2} \sum_{a \in S} \sum_{b \in S} 1_{\{a \in S_A \text{ and } b \in S_B\}}} \xrightarrow{P} \frac{|E_{A,B}|}{|A||B|} = w(A, B) \\ &\text{by LLN and Slutsky's Theorem (p.3).} \end{aligned}$$

$$\begin{aligned} \text{Eq.(9): } \hat{w}(A, B) &= \frac{\sum_{a \in S_A} |E_{a,B}| + \sum_{b \in S_B} |E_{b,A}|}{|S_A| \cdot |\hat{B}| + |S_B| \cdot |\hat{A}|} = \\ &= \frac{\frac{N}{n} \sum_{a \in S_A} |E_{a,B}| + \frac{N}{n} \sum_{b \in S_B} |E_{b,A}|}{\frac{N}{n} |S_A| \cdot |\hat{B}| + \frac{N}{n} |S_B| \cdot |\hat{A}|} \xrightarrow{P} \frac{|E_{A,B}| + |E_{B,A}|}{|A||B| + |A||B|} = \\ &w(A, B) \text{ by the LLN on numerator and denominator and then by five applications of Slutsky's Theorem (p.1, p.2, and p.3).} \end{aligned}$$

Non-uniform sampling estimators

$$\text{Eq.(10): } \hat{x}_{\text{tot}} = \frac{1}{n} \sum_{v \in S} \frac{x(v)}{\pi(v)} \text{ is shown to be consistent in [25].}$$

$$\text{Eq.(11): } |\hat{A}| = N \cdot \frac{\frac{1}{n} w_1(S_A)}{\frac{1}{n} w_1(S)} \xrightarrow{P} |A| \text{ by the consistency of the HH estimator in the numerator and denominator and then by Slutsky's Theorem (p.3).}$$

$$\text{Eq.(14): } \hat{k}_V = \frac{\frac{1}{n} \sum_{v \in S} \frac{\deg(v)}{w(v)}}{\frac{1}{n} w_1(S)} \xrightarrow{P} k_V \text{ and}$$

$$\hat{k}_A = \frac{\frac{1}{n} \sum_{v \in S_A} \frac{\deg(v)}{w(v)}}{\frac{1}{n} w_1(S_A)} \xrightarrow{P} k_A \text{ by the consistency of the HH estimator and Slutsky's Theorem (p.3).}$$

$$\begin{aligned} \text{Eq.(13): } \hat{f}_A^{\text{vol}} &= \frac{\frac{1}{n} \sum_{s \in S} \left(\frac{1}{w(s)} \sum_{v \in \mathcal{N}(s)} 1_{\{v \in A\}} \right)}{\frac{1}{n} \sum_{s \in S} \frac{\deg(s)}{w(s)}} \xrightarrow{P} f_A^{\text{vol}} \\ &\text{by the consistency of the HH estimator and Slutsky's Theorem (p.3).} \end{aligned}$$

$$\text{Eq.(12): } |\hat{A}| = N \cdot \hat{f}_A^{\text{vol}} \cdot \frac{\hat{k}_V}{\hat{k}_A} \xrightarrow{P} |A| \text{ by the consistency of the estimators and Slutsky's Theorem (p.2 and p.3).}$$

$$\begin{aligned} \text{Eq.(15): } \hat{w}(A, B) &= \frac{\frac{1}{n^2} \sum_{a \in S_A} \sum_{b \in S_B} \frac{1_{\{a, b\} \in E}}{w(a) \cdot w(b)}}{\frac{1}{n^2} w_1(S_A) \cdot w_1(S_B)} \\ &\xrightarrow{P} \frac{|E_{A,B}|}{|A||B|} = w(A, B) \text{ by the consistency of the HH estimator and Slutsky's Theorem (p.2 and p.3).} \end{aligned}$$

$$\begin{aligned} \text{Eq.(16): } \hat{w}(A, B) &= \frac{\frac{1}{n} \sum_{a \in S_A} \frac{|E_{a,B}|}{w(a)} + \frac{1}{n} \sum_{b \in S_B} \frac{|E_{b,A}|}{w(b)}}{\frac{1}{n} w_1(S_A) \cdot |\hat{B}| + \frac{1}{n} w_1(S_B) \cdot |\hat{A}|} = \\ &\xrightarrow{P} \frac{|E_{A,B}| + |E_{B,A}|}{|A||B| + |A||B|} = w(A, B) \text{ by the consistency of HH estimators in the numerator and denominator and then by five applications of Slutsky's Theorem (p.1, p.2 and p.3).} \end{aligned}$$

A note on dependent samples

These results continue hold in the case of dependent (correlated) samples, such as RW, under the condition that these samples converge asymptotically to UIS or WIS limits. This follows from the ergodic theorem, which provides a corresponding LLN for convergent Markov Chains. For more technical details on the LLN in the context of dependent samples see [52].